

<https://helda.helsinki.fi>

High-resolution population-specific recombination rates and their effect on phasing and genotype imputation

Hassan, Shabbeer

2021-04

Hassan , S , Surakka , I , Taskinen , M-R , Salomaa , V , Palotie , A , Wessman , M ,
Tukiainen , T , Pirinen , M , Palta , P & Ripatti , S 2021 , ' High-resolution population-specific
recombination rates and their effect on phasing and genotype imputation ' , European
pöJournal of Human Genetics , vol. 29 , no. 4 , pp. 615 624 . <https://doi.org/10.1038/s41431-020-00768-8>

<http://hdl.handle.net/10138/332029>

<https://doi.org/10.1038/s41431-020-00768-8>

cc_by

publishedVersion

Downloaded from Helda, University of Helsinki institutional repository.

This is an electronic reprint of the original article.

This reprint may differ from the original in pagination and typographic detail.

Please cite the original version.



ARTICLE

High-resolution population-specific recombination rates and their effect on phasing and genotype imputation

Shabbeer Hassan¹ · Ida Surakka¹ · Marja-Riitta Taskinen² · Veikko Salomaa³ · Aarno Palotie^{1,4,5} ·
Maija Wessman¹ · Taru Tukiainen¹ · Matti Pirinen^{1,6,7} · Priit Palta^{1,8} · Samuli Ripatti^{1,5,6}

Received: 27 May 2020 / Revised: 1 October 2020 / Accepted: 20 October 2020 / Published online: 28 November 2020
© The Author(s) 2020. This article is published with open access

Abstract

Previous research has shown that using population-specific reference panels has a significant effect on downstream population genomic analyses like haplotype phasing, genotype imputation, and association, especially in the context of population isolates. Here, we developed a high-resolution recombination rate mapping at 10 and 50 kb scale using high-coverage (20–30×) whole-genome sequenced data of 55 family trios from Finland and compared it to recombination rates of non-Finnish Europeans (NFE). We tested the downstream effects of the population-specific recombination rates in statistical phasing and genotype imputation in Finns as compared to the same analyses performed by using the NFE-based recombination rates. We found that Finnish recombination rates have a moderately high correlation (Spearman's $\rho = 0.67$ – 0.79) with NFE, although on average (across all autosomal chromosomes), Finnish rates (2.268 ± 0.4209 cM/Mb) are 12–14% lower than NFE (2.641 ± 0.5032 cM/Mb). Finnish recombination map was found to have no significant effect in haplotype phasing accuracy (switch error rates ~2%) and average imputation concordance rates (97–98% for common, 92–96% for low frequency and 78–90% for rare variants). Our results suggest that haplotype phasing and genotype imputation mostly depend on population-specific contexts like appropriate reference panels and their sample size, but not on population-specific recombination maps. Even though recombination rate estimates had some differences between the Finnish and NFE populations, haplotyping and imputation had not been noticeably affected by the recombination map used. Therefore, the currently available HapMap recombination maps seem robust for population-specific phasing and imputation pipelines, even in the context of relatively isolated populations like Finland.

Supplementary information The online version of this article (<https://doi.org/10.1038/s41431-020-00768-8>) contains supplementary material, which is available to authorized users.

✉ Samuli Ripatti
samuli.ripatti@helsinki.fi

- ¹ Institute for Molecular Medicine Finland, FIMM, HiLIFE, University of Helsinki, Helsinki, Finland
- ² Clinical and molecular metabolism, Research program unit, University of Helsinki, Helsinki, Finland
- ³ Finnish Institute for Health and Welfare, Helsinki, Finland
- ⁴ Massachusetts General Hospital & Harvard Medical School, Boston, MA, USA
- ⁵ Broad Institute of the Massachusetts Institute of Technology and Harvard University, Cambridge, MA, USA
- ⁶ Department of Public Health, Faculty of Medicine, Clinicum, University of Helsinki, Helsinki, Finland
- ⁷ Department of Mathematics and Statistics, University of Helsinki, Helsinki, Finland
- ⁸ Estonian Genome Center, Institute of Genomics, University of Tartu, Tartu, Estonia

Introduction

Recombination is not uniform across the human genome with large areas having lower recombination rates, so-called ‘coldspots’, which are then interspersed by shorter regions marked by a high recombinational activity called ‘hotspots’ [1]. With long chunks of human genome existing in high linkage disequilibrium, LD [2], and organised in the form of ‘haplotype blocks’, the ‘coldspots’ tend to coincide with such regions of high LD [3].

Direct estimation methods of recombination are quite time-consuming, and evidence has suggested that they do not easily scale up to genome-wide, fine-scale recombinational variation estimation [4]. A less time-consuming but computationally intensive alternative is to use the LD patterns surrounding the SNPs [5]. Such methods have been used in the past decade or so, to create fine-scale recombination maps [6]. Besides the International HapMap project that focused on capturing common variants and haplotypes in diverse populations, international whole-

genome sequencing (WGS)-based collaborations like the 1000 Genomes Project, provided genetic variation data for 20 worldwide populations [7]. This led to further refinement of the recombination maps coupled with methodological advances of using coalescent methods for recombination rate [8, 9].

With the rise of international collaborative projects, it was realised that founder populations can often have very unique LD patterns [10], subsequently also displaying unique increased genetics-driven health risks [11], suggesting that population-specific reference datasets should be used to leverage the LD patterns to better detect disease variants in the downstream genetic analysis [12]. Genomic analysis methods like haplotype phasing and imputing genotypes require recombination maps and other population genetic parameters as input to obtain optimal results [13–16].

In this study, we set to test this by (1) estimating recombination rates along the genome in Finnish population using ~55 families of whole-genome sequenced (20–30×) Finns, (2) comparing these rates to some other European populations, and (3) comparing the effect of using Finnish recombination rate estimates and cosmopolitan estimates in phasing and imputation errors in Finnish samples.

Materials and methods

Datasets used

Finnish migraine families collection

Whole-genome sequenced trios ($n = 55$) consisting of the parent-offspring combination were drawn from a large Finnish migraine families collection consisting of 1589 families totalling 8319 individuals [17]. The trios were used for the recombination map construction using LDHAT version 2. The families were collected over 25 years from various headache clinics in Finland (Helsinki, Turku, Jyväskylä, Tampere, Kemi, and Kuopio) and via advertisements in the national migraine patient organisation web page (<https://migreeni.org/>). The families consist of different pedigree sizes from small to large (1–5+ individuals). Of the 8319 individuals, 5317 have a confirmed migraine diagnosis based on the third edition of the established International Classification for Headache Disorders (ICHD-3) criteria [18].

EUFAM cohort

To check the phasing accuracy of our Finnish recombination map, we used an independently sourced 49 trios from the European Multicenter Study on Familial Dyslipidemias

in Patients with Premature Coronary Heart Disease (EUFAM). Finnish familial combined hyperlipidemia families were identified from patients initially admitted to hospitals with premature cardiovascular heart disease diagnosis who also had elevated levels of total cholesterol, triglycerides (TG), or both in the ≥90th Finnish population percentile. Those families who had at least one additional first-degree relative also affected with hyperlipidemia were also included in the study apart from individuals with elevated levels of TG [19–21].

FINRISK cohort

The imputation accuracy of the Finnish and previously published HapMap based recombination maps [8, 9] was subsequently tested on an independent FINRISK CoreExome chip dataset consisting of 10,481 individuals derived from the national-level FINRISK cohort. Primarily, it comprises of respondents of representative, cross-sectional population surveys that are conducted once every 5 years since 1972 to get a national assessment of various risk factors of chronic diseases and other health behaviours among the working-age population drawn from 3 to 4 major cities in Finland [22].

FINNISH genomic reference panel cohort

The whole-genome sequenced samples used were obtained from PCR-free methods and PCR-amplified methods, which was followed by sequencing on an Illumina HiSeq X platform with a mean depth of ~30×. The obtained reads were then aligned to the GRCh37 (hg19) human reference genome assembly using BWA-MEM. Best practice guidelines from Genome Analysis Toolkit (GATK) were used to process the BAM files and variant calling. Several criteria were used in this stage for sample exclusion: relatedness (identity-by-descent (IBD) > 0.1), sex mismatches, among several others. Furthermore, samples were filtered based on other criteria such as non-reference variants, singletons, heterozygous/homozygous variants ratio, insertion/deletion ratio for novel indels, insertion/deletion ratio for indels observed in dbSNP, and transition/transversion ratio.

After this stage, some exclusion criteria were applied to set some variants as missing: $GQ < 20$, phred-scaled genotype likelihood of reference allele < 20 for heterozygous and homozygous variant calls, and allele balance < 0.2 or > 0.8 for heterozygous calls. A truth sensitivity percentage threshold of 99.8% for SNVs and of 99.9% for indels was used based on the GATK Variant Quality Score Recalibration to filter variants with, quality by depth < 2 for SNVs and < 3 for indels, call rate < 90%, and Hardy–Weinberg equilibrium p -value < 1×10^{-9} . Some other variants like

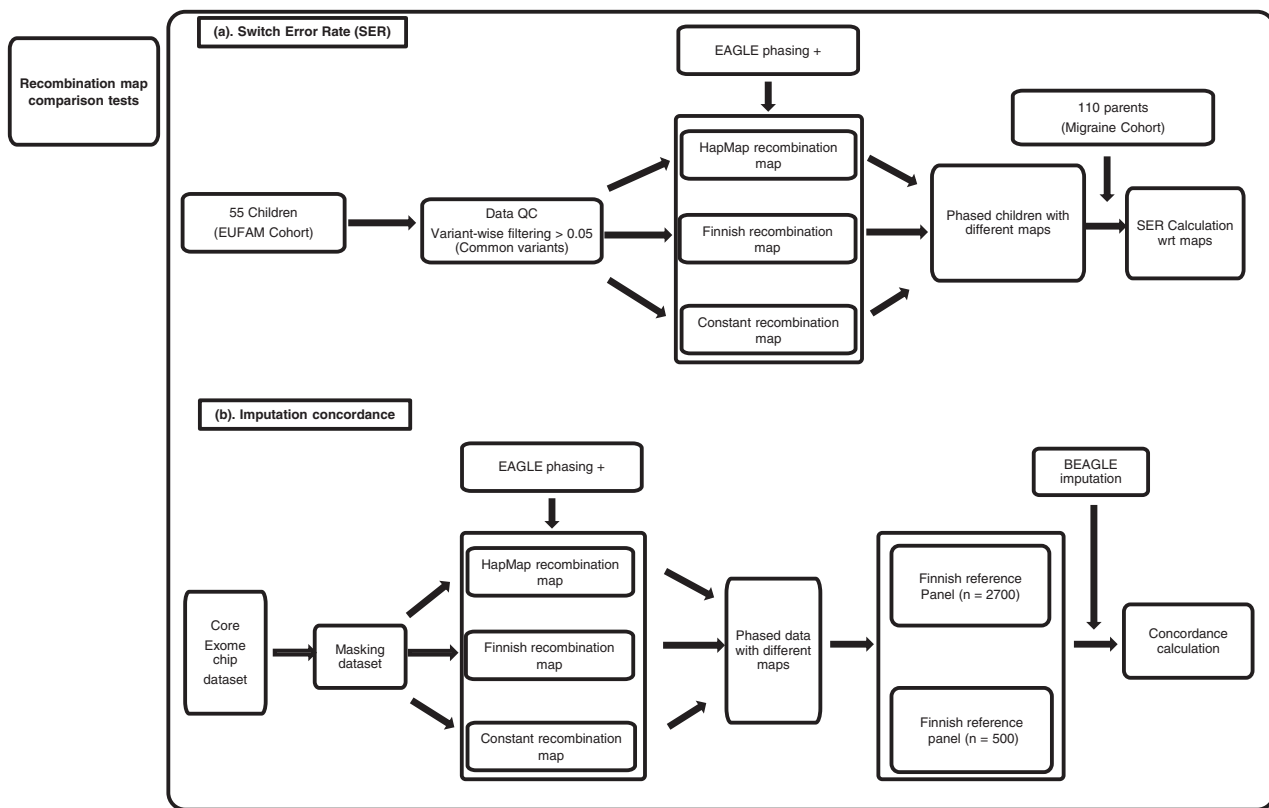


Fig. 1 Flowchart. Overview of the analyses and comparisons performed.

monomorphic, multi-allelic, and low-complexity regions [23] were further excluded.

The final reference dataset (SISu v2) used in this study for imputation consisted of high-coverage (20–30×) whole-genome sequence-based reference panel of 2690 individuals from the SISu project (Sequencing Initiative Suomi, <http://www.sisuproject.fi/>, [24]). Here, SISu v2 refers to in-house version 2 (with 2690 individuals) of our SISu imputation reference panel.

Recombination map construction

Coalescent-based fine-scale recombination map construction [8] is greatly eased by using trios which provide more accurate haplotype phasing resolution [25]. Hence, we used trio data ($n = 55$, 110 independent parents) from the Finnish Migraine Families Cohort described above. These were filtered primarily using VCFtools [26] and custom R scripts. Firstly, sites were thinned with within 15 bp of each other such that only one site remained followed by a filtering step of removing variants with a minor allele frequency of <5% [27]. The resultant data were then phased using a family-aware method of SHAPEIT [28] using the standard HapMap recombination map [8, 9], which was then split into segments of ~10,000 SNPs with a 1000 SNP overhang on

each side of the segments. LDhat version 2 was run for 10^7 iterations with a block penalty of 5, every 5000 iterations of them of which the first 10% observations were discarded [8, 29]. The CEU based maps, used here for comparison, were obtained similarly using LDhat [29].

However, LDhat is computationally intensive, and calculations suggest that the 1000 Genomes OMNI dataset [30] would be too much computationally intensive to complete [31], hence limiting the maximum number of haplotypes which could be used.

To overcome this and make the recombination map independent of the underlying methodology, we used a machine learning method implemented in FastEPRR [31, 32]. It supports the use of larger sample sizes, than LDhat and the recombination estimates for sample sizes >50, yields smaller variance than LDhat-based estimates [31]. The method was then applied to each autosome with overlapping sliding windows (i.e., window size, 50 kb and step length, 25 kb) under default settings for diploid organisms. As seen in [31] both methods produce similar estimates, with variance of the estimate of mean being different.

The output of LDhat and FastEPRR is in terms of population-recombination rate (p) and to convert them into per-generational rate (r) used in phasing/imputation

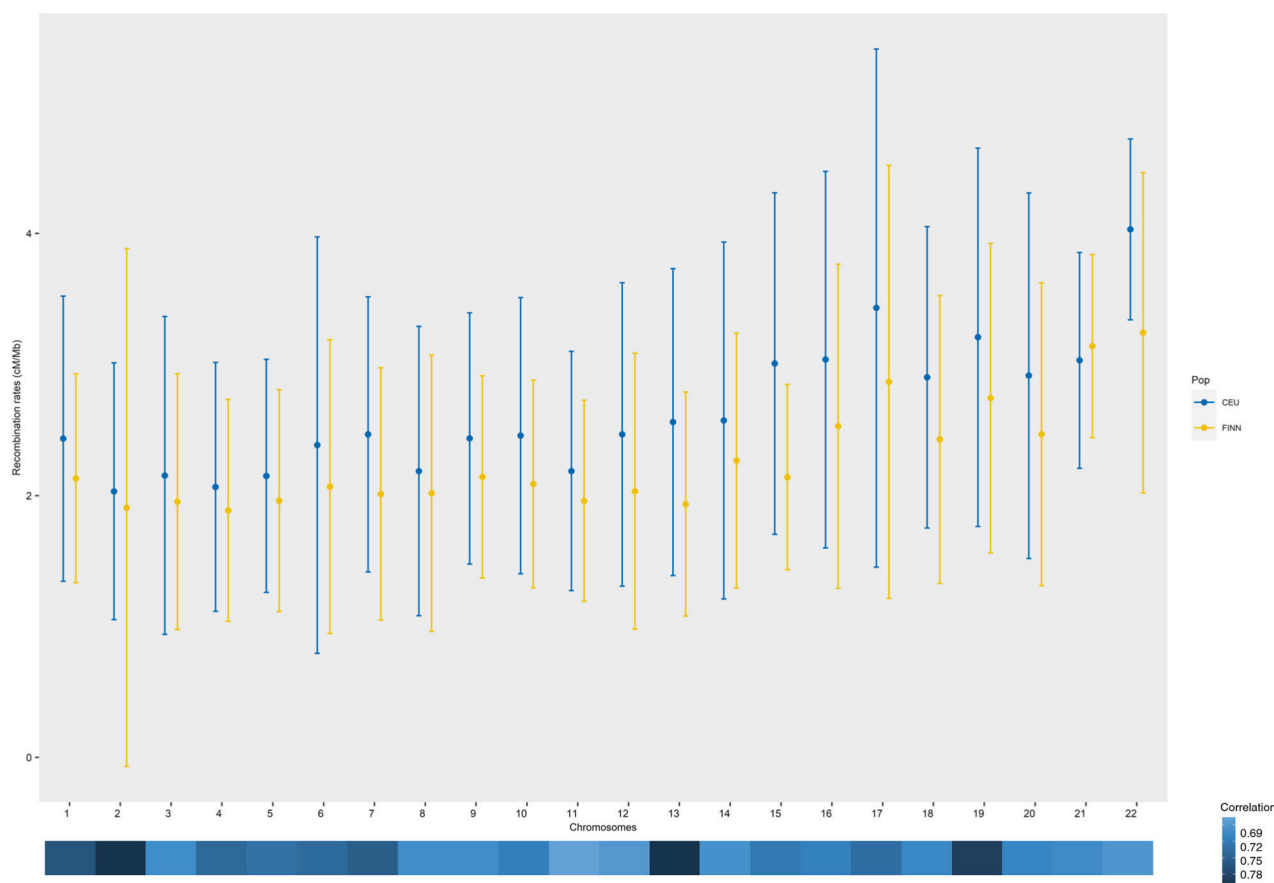


Fig. 2 Average (\pm standard deviation) recombination rates of Finnish v/s CEU per autosome measured in cM/Mb and correlation between Finnish and CEU recombination rates across all chromosomes. The comparisons are made for similar physical positions.

algorithms we used optimal effective population size values derived from our testing (as explained in the Supplementary Text). The estimates from LDHAT and FastEPRR were then averaged, to obtain a new combined estimate with the lowest variance amongst all the three [31].

Phasing and imputation accuracy

To test whether the usage of different recombination maps affects the efficiency of haplotype phasing and imputation, we used the aforesaid Finnish genotype data to evaluate: (i) switch error rates (SER) across all chromosomes and (ii) imputation concordance rates for chromosome 20.

Phasing accuracy

The gold standard method to estimate haplotype phasing accuracy is to count the number of switches (or recombination events) needed between the computationally phased dataset and the true haplotypes [33]. The number of such switches divided by the number of all possible switches is called SER.

For testing the influence of recombination maps on phasing accuracy, we used three different recombination maps: HapMap, fine-scale Finnish recombination map, and a constant background recombination rate (1 cM/Mb), to phase the 55 offspring haplotypes without using any reference dataset. To check whether reference panels used during haplotype phasing made any impact on the SER, we used the Finnish SISU based reference ($n = 2690$), to check whether the size of the reference panel made any impact on the results in phasing the offspring's haplotypes (Fig. 1).

The SER in the offspring's phased haplotypes were then calculated by determining the true offspring haplotypes using data from the parents (98 individuals) with a custom script [34].

Imputation accuracy

Imputation concordance was used as the metric for calculating the imputation accuracy. For this, we randomly masked FINRISK CoreExome chip data consisting of 10,480 individuals [22] from chromosome 20. To test the

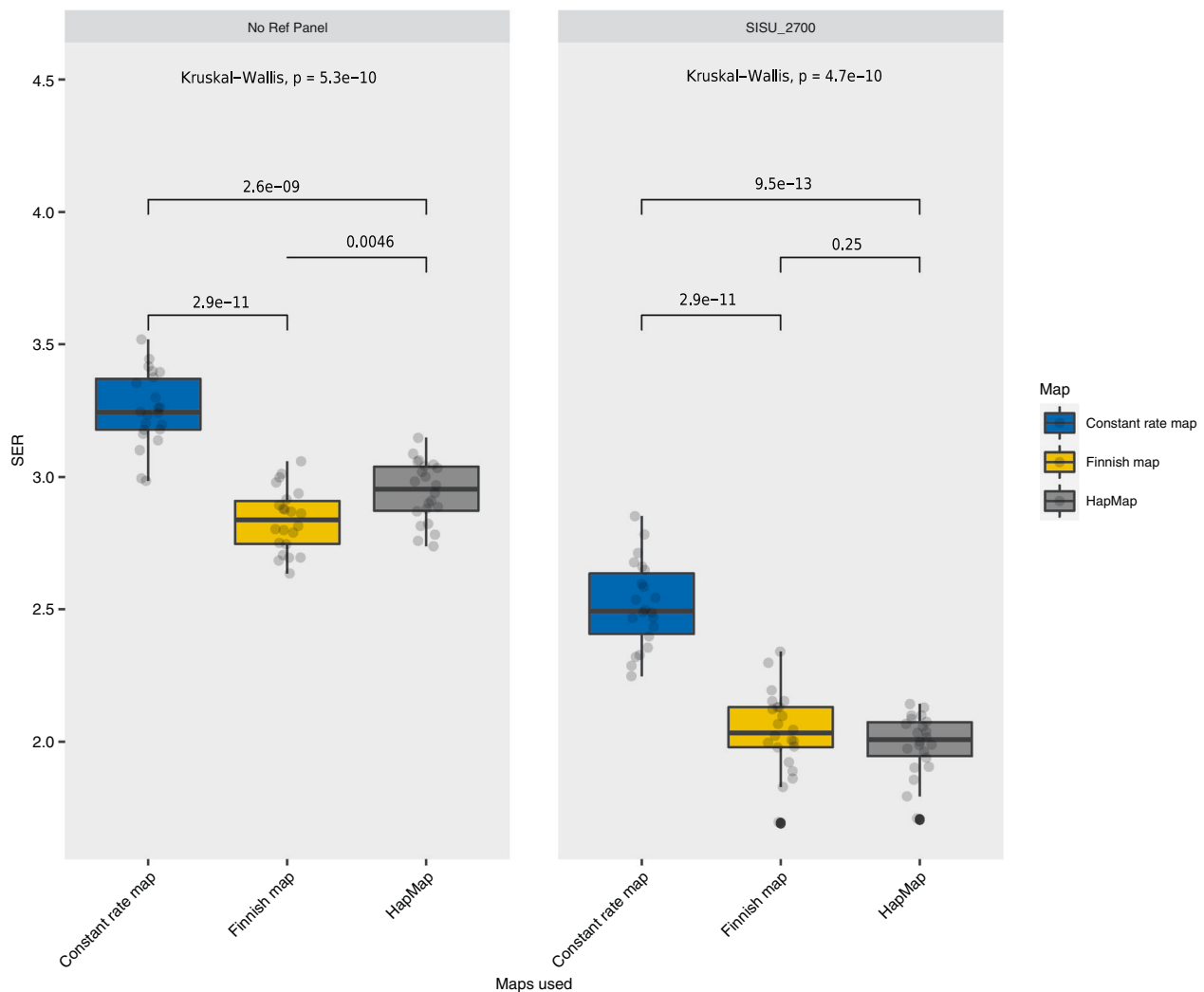


Fig. 3 Statistical comparison of switch error rates across all autosomes calculated for all children in the trios using different recombination maps with respect to different reference panel

role of reference panel size in influencing the imputation accuracy in conjunction with varying the population genetics parameters, we imputed the masked dataset with BEAGLE [15] using the Finnish reference panel ($n = 2690$). The concordance was then calculated between the imputed genotypes and the original masked variants. Masking was done by randomly removing ~10% of variants from the chip dataset.

The influence of recombination maps on imputation accuracy was checked by calculating the concordance values between imputed and original variants, using the Finnish reference panel in various combinations of recombination maps (constant rate, HapMap, Finnish map) during the imputation (Fig. 1). Constant rate map used here consisted of a constant recombination rate of 1 cM/Mb, used as a control condition for testing the statistical differences.

conditions (absent or present). The p -values are shown at the top of each panel from Kruskal Wallis ANOVA testing between panel groups and ones between boxplots for within-group comparisons.

Results

Finnish recombination map and its comparison to the HapMap recombination map

The primary aim of our study was to derive a high-resolution genetic recombination map for Finland and use it for comparative tests in commonly used analyses like haplotype phasing and imputation. To derive a population-specific Finnish recombination map, we used the high-coverage WGS data and an average of different estimation methods (LDHat and FastEPFR). We used the N_e value of 10,000 derived from our extensive testing of different N_e values (See Supplementary Text) to get the per-generation recombination rates. The average recombination rates of Finnish population isolate depicted 12–14% lower values (autosome-wide average 2.268 ± 0.4209 cM/Mb) for all

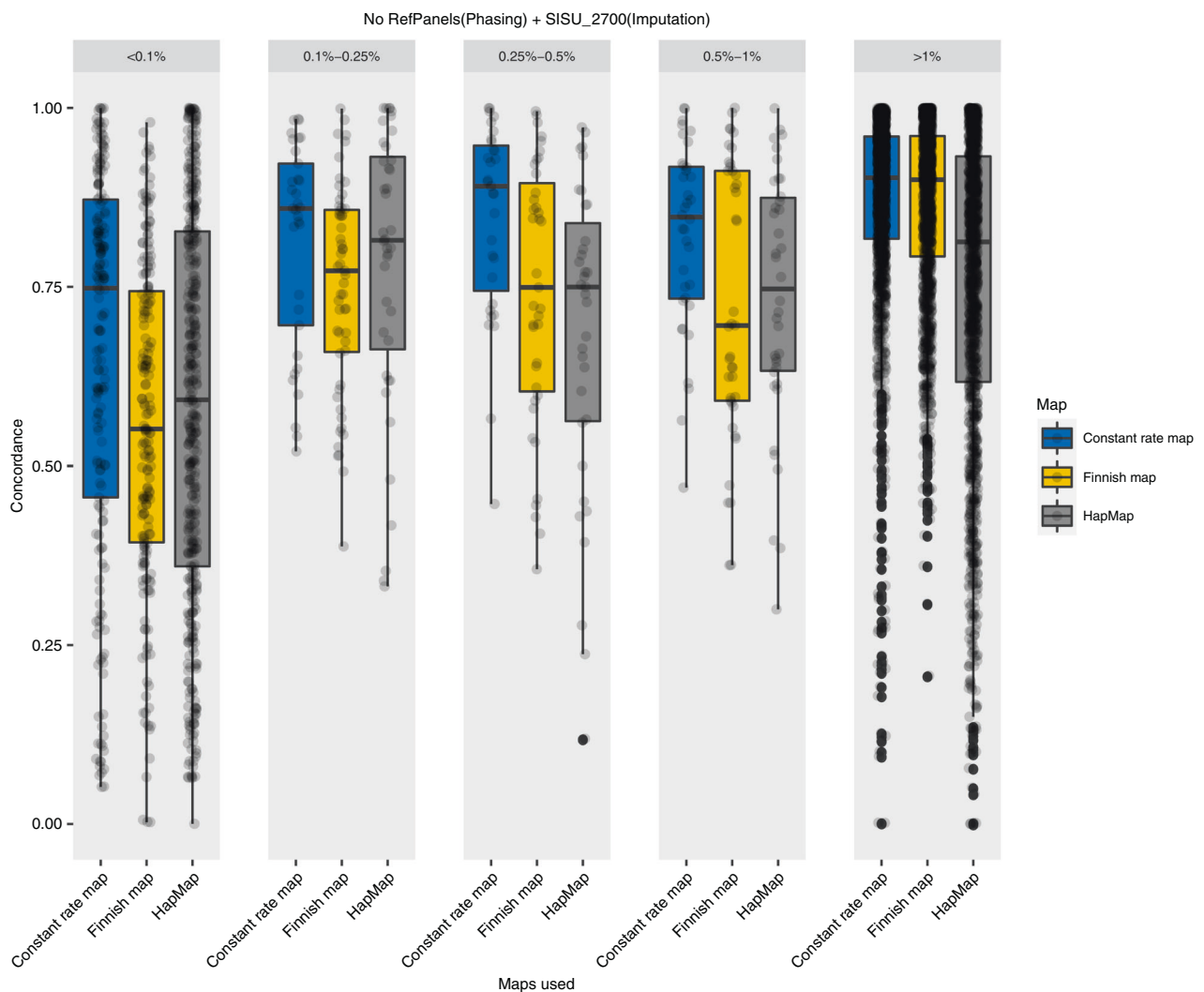


Fig. 4 Imputation concordance (NO reference panels). Comparison across different minor allele frequency (MAF) groups for a range of different recombination map combinations phased with no reference panels.

chromosomes compared to CEU based maps (2.641 ± 0.5032 cM/Mb) (Fig. 2).

These differences in average recombination rates are reflected in the correlation values across all chromosomes (Spearman's $\rho \sim 0.67$ – 0.79) between the developed Finnish map and HapMap based one (Fig. 2). We also present a direct comparison between the two maps, of the recombination rates at 5 Mb scales, which presents a similar visual pattern of rates across the genome (Supplementary Fig. 1).

Effects of the population-specific recombinations map on haplotype phasing

Variation in population-specific recombination maps (and effective population sizes) can affect the downstream genomic analyses like haplotype phasing and imputation.

We tested the Finnish map, HapMap map, and a constant recombination rate map (1 cM/Mb) to understand the effects of population-specific maps on downstream genomic analyses. The phasing accuracy was tested under two different conditions: using no additional reference panel and using a population-specific SISu v2 reference panel ($n = 2690$) in phasing. We observed that, on average, SER ranged between 1.8 and 3.7% (Supplementary Fig. 2) across the different chromosomes and recombination maps. We found statistically significant differences within both no reference panel and the Finnish reference panel results (Kruskal Wallis, p -value = $5.3e-10$ and $4.7e-10$, respectively; Fig. 3). The constant recombination map (1 cM/Mb) had significantly higher SER values when compared to the Finnish map or the HapMap map (Fig. 3) both when no reference panels were used (p -value = $2.9e-11$ and $2.6e-09$, respectively) and when the Finnish reference panel

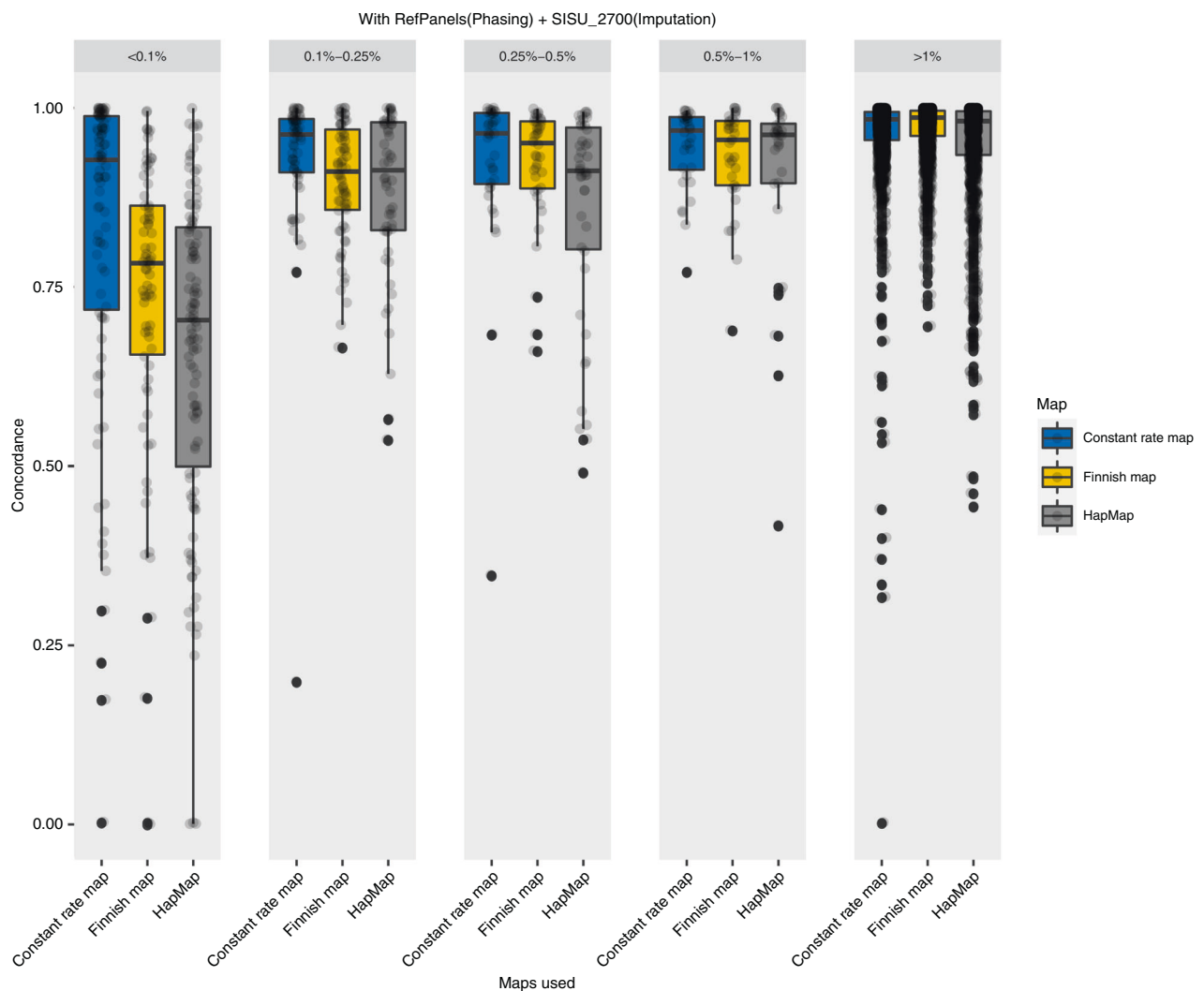


Fig. 5 Imputation concordance (WITH reference panels). Comparison across different minor allele frequency (MAF) groups for a range of different recombination map combinations phased with reference panels.

was used (p -value = $2.9\text{e}-11$ and $9.5\text{e}-13$, respectively). The choice of recombination maps mattered more when no reference panel was used (p -value = 0.0046), however when using the Finnish reference panel, the difference in SER was statistically insignificant (p -value = 0.25).

Effects of the population-specific recombinations map on genotype imputation

Imputation accuracy was similarly tested using the reference panel under three different recombination map settings. We observed that when the imputation target dataset was phased and imputed using the Finnish reference panel ($n = 2690$) irrespective of the population-specific recombination maps, it had a high imputation accuracy (overall concordance rate ~98%, Fig. 4) across MAF bins ($>0.1\%$). Though some differences in concordance rates are seen in

for rare variants (MAF $<0.1\%$). The concordance rate was lower when the test dataset was phased without reference panels (concordance rate 72–77%, Fig. 5).

Discussion

Population isolates like Finland, have had a divergent demographic history as compared to the outbred Non-Finnish European populations, with lower migration rates, more fluctuating population sizes and higher incidences of bottleneck events and founder effects [35, 36]. This unique demographic history then affects different population genetic parameters, like recombination rates [37]. It has been shown previously that using population-specific genomic reference panels augmented the accuracy of imputation accuracy leading to better mapping of diseases

specific variants in GWAS [12]. Since recombination rates (in the form of recombination maps), features in much of the downstream genomic analyses' methods like imputation and haplotype phasing [15, 34], we wanted to study their effect on downstream analyses.

Firstly, we characterised the Finnish recombination map using high-coverage (~30×) WGS samples from large SISu v2 reference panel ($n = 2690$). Previously used recombination maps hail from the HapMap and 1000 Genomes projects which used sparse genotypic datasets or low-depth sequencing samples. This is a first attempt in creating a recombination map for Finland using population-specific WGS samples. We used two different methods in estimating the recombination rates, to achieve accurate estimates with lower variance [29, 31]. We also estimated effective population sizes using IBD based methods [15] for both Finnish and CEU based datasets. The obtained recombination map was then used to test their role and importance in two selected downstream genomic analyses—haplotype phasing and imputation concordance. Since the recombination rate determination requires effective population size estimates, we also tested the role of varying effective population size on these two analyses (See Supplementary Text). The extensive testing of N_e yielded the estimate of 10,000 originally derived theoretically [38] and most used commonly for humans fits quite rightly for the recombination map.

The Finnish recombinational landscape when compared to the HapMap based map, showed, on average, a high degree of correlation across scales (10, 50 kb, and 5 Mb), however, on average, Finnish recombination rates across chromosomes were found to be lower. Such moderate to high correlations (Fig. 2) and similar recombinational landscape (Supplementary Fig. 1) could be due to high sharing of recombinations in individuals from closely related populations. The degree of dissimilarity in the population-level differences between Finnish and mainland Europeans in terms of recombination rates could be due to population-specific demographic processes like founder effects, bottleneck events and migration [39], and other biological processes directly affecting the recombination rate [40]. The broad similarity in terms of correlational structure observed here reflects a shared ancestral origin of Finns and mainland Europeans [41]. Other studies on population isolates like Iceland [9] have previously found a high degree of correlation with CEU based maps, albeit with substantial differences as seen here. Previous studies [42] have additionally explored the relationship between recombination rate differences between populations and allele frequency differences, with evidence suggesting that the differences between rates show the selection impact in the past 100,000 years since the out-of-Africa movement of humans.

As seen in previous studies, much of the downstream genomic analyses like getting more refined GWAS hits or, accurate copy number variants imputation, can be highly improved with the addition/use of population-specific datasets [12]. To test this in the context of population-specific recombination maps, we used them to test the haplotype phasing and imputation accuracy and observed that despite substantial differences in the effective population sizes between populations, it did not affect the tested metrics. One possible explanation for the insignificant effect seen here is that the role of parameters like effective population size and recombination maps is to scale over the haplotypes for efficient coverage of the whole genome. However, when sufficiently large, population-specific genomic reference panels are available with tens of thousands of haplotypic combinations, such scaling over for specific populations, does not yield in substantial improvements. As we showed here, the selection of reference panel size could play an important role in the downstream genomic analyses and for most cases, the current practice of using the standard HapMap recombination map can be reasonably used. Another point of interest here is that the use of different N_e parameters during phasing/imputation might be redundant as we observed no change in the accuracy of our estimates on varying the N_e parameters. Similarly, when using population-specific recombination maps in haplotype phasing or genotype imputation, we did not find any tangible benefits in using them over the current standard maps based on the HapMap data. On the other hand, population-specific recombination rates could be quite important for other population genetic processes and their estimators.

One of the main limitations of our study comes from the LD-based methods used in the estimation of population-specific recombination rates. Several studies have shown that different demographic history between populations, can askew/bias the estimates derived therein [43–45]. Several recent methods do exist which overcome this limitation history [46, 47], and using them might lead to slightly unbiased different estimates. Hence, even though we do find a difference in autosomal-wide recombination rates between Finnish and NFE, these could also be interpreted as due to the inherent bias in the estimation of population-recombination rate (ρ). Our study suggests a couple of important points for future studies: (a) varying effective population size for downstream genomic analyses, such as phasing and imputation, might have a relatively small impact, and it might be better to use the default option of the particular software; (b) when available, it is beneficial to use a population-specific genomic reference panel as they increase the accuracy; (c) HapMap can be used for current downstream genomic analyses like haplotype phasing or genotype imputation in European-based populations. And,

if need be, can be substituted for using population-specific maps, as the accuracy rates are quite similar to the population-based maps.

Though the sample used here is from a disease cohort but is nevertheless representative of Finland's population and hence provides a reasonable recombination rate estimates. On the other hand, our reliance on disease cohorts could lead to minor variation in the resultant recombination.

To conclude, even though our estimates of recombination rates had some differences between the Finnish and NFE populations, and the effective population size may vary between these two populations, we did not observe that the downstream analyses of haplotyping and imputation had been noticeably affected by the recombination map or the values of effective population size used.

Acknowledgements We would like to thank Sari Kivikko and Huei-Yi Shen for management assistance. The FINRISK analyses were conducted using the THL biobank permission for project BB2015_55.1. We thank all study participants for their generous participation in the FINRISK study.

Funding This work was financially supported by the Academy of Finland (251217 and 255847 to SR). SR was further supported by the Academy of Finland, Center of Excellence for Complex Disease Genetics, the Finnish Foundation for Cardiovascular Research, Biocentrum Helsinki, and the Sigrid Jusélius Foundation. SH was supported by FIMM-EMBL PhD program doctoral funding and IS by Academy of Finland Postdoctoral Fellowship (298149). VS was supported by the Finnish Foundation for Cardiovascular Research. TT was supported by Academy of Finland grant number 315589.

Compliance with ethical standards

Conflict of interest VS has received honoraria from Novo Nordisk and Sanofi for consulting and has an ongoing research collaboration with Bayer Ltd (all unrelated to the present study).

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- Baudat F, Buard J, Grey C, Fledel-Alon A, Ober C, Przeworski M, et al. PRDM9 is a major determinant of meiotic recombination hotspots in humans and mice. *Science*. 2009;327:836–40.
- Daly MJ, Rioux JD, Schaffner SF, Hudson TJ, Lander ES. High-resolution haplotype structure in the human genome. *Nat Genet*. 2001;29:229–32.
- Hudson RR, Kaplan NL. Statistical properties of the number of recombination events in the history of a sample of DNA sequences. *Genetics*. 1985;111:147–64.
- Chan AH, Jenkins PA, Song YS. Genome-wide fine-scale recombination rate variation in *Drosophila melanogaster*. *PLoS Genet*. 2012;8:e1003090.
- McVean GA, Myers SR, Hunt S, Deloukas P, Bentley DR, Donnelly P. The fine-scale structure of recombination rate variation in the human genome. *Science*. 2004;304:581–4.
- Myers S, Bottolo L, Freeman C, McVean G, Donnelly P. A fine-scale map of recombination rates and hotspots across the human genome. *Science*. 2005;310:321–4.
- Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel JO, et al. A global reference for human genetic variation. *Nature*. 2015;526:68–74.
- Auton A, McVean G. Recombination rate estimation in the presence of hotspots. *Genome Res*. 2007;17:1219–27.
- Kong A, Thorleifsson G, Gudbjartsson DF, Masson G, Sigurdsson A, Jonasdottir A, et al. Fine-scale recombination rate differences between sexes, populations and individuals. *Nature*. 2010;467:1099–103.
- Service S, DeYoung J, Karayiorgou M, Roos JL, Pretorius H, Bedoya G, et al. Magnitude and distribution of linkage disequilibrium in population isolates and implications for genome-wide association studies. *Nat Genet*. 2006;38:556–60.
- Peltonen L, Jalanko A, Varilo T. Molecular genetics of the Finnish disease heritage. *Hum Mol Genet*. 1999;8:1913–23.
- Surakka I, Kristiansson K, Anttila V, Inouye M, Barnes C, Moutsianas L, et al. Founder population-specific HapMap panel increases power in GWA studies through improved imputation accuracy and CNV tagging. *Genome Res*. 2010;20:1344–51.
- Tewhey R, Bansal V, Torkamani A, Topol EJ, Schork NJ. The importance of phase information for human genomics. *Nat Rev Genet*. 2011;12:215–23.
- Browning SR, Browning BL. Haplotype phasing: existing methods and new developments. *Nat Rev Genet*. 2011;12:703–14.
- Browning BL, Browning SR. Genotype imputation with millions of reference samples. *Am J Hum Genet*. 2016;98:116–26.
- Delaneau O, Zagury JF, Marchini J. Improved whole-chromosome phasing for disease and population genetic studies. *Nat Methods*. 2013;10:5–6.
- Gormley P, Kurki MI, Hiekkala ME, Veerapen K, Häppölä P, Mitchell AA, et al. Common variant burden contributes to the familial aggregation of migraine in 1,589 families. *Neuron*. 2018;98:743–e4.
- Headache Classification Committee of the International Headache Society. The International Classification of Headache Disorders, 3rd edn. (beta version). *Cephalalgia* 2013;33:629–808.
- Borodulin K, Vartiainen E, Peltonen M, Jousilahti P, Juolevi A, Laatikainen T, et al. Forty-year trends in cardiovascular risk factors in Finland. *Eur J Public Health*. 2015;25:539–46.
- Porkka KV, Nuotio I, Pajukanta P, Ehnholm C, Suurinkeroinen L, Syväne M, et al. Phenotype expression in familial combined hyperlipidemia. *Atherosclerosis*. 1997;133:245–53.
- Ripatti P, Rämö JT, Söderlund S, Surakka I, Matikainen N, Pirinen M, et al. The contribution of GWAS loci in familial dyslipidemias. *PLOS Genet*. 2016;12:e1006078.
- Vartiainen E, Laatikainen T, Peltonen M, Juolevi A, Mannisto S, Sundvall J, et al. Thirty-five-year trends in cardiovascular risk factors in Finland. *Int J Epidemiol*. 2009;39:504–18.
- Li H. Toward better understanding of artifacts in variant calling from high-coverage samples. *Bioinformatics*. 2014;30:2843–51.

24. Kals M, Nikopoulou T, Läll K, Pärn K, Sikka TT, Suvisaari J, et al. Advantages of genotype imputation with ethnically matched reference panel for rare variant association analyses bioRxiv. 579201. <https://doi.org/10.1101/579201>.
25. Roach JC, Glusman G, Hubley R, Montsaroff SZ, Holloway AK, Mauldin DE, et al. Chromosomal haplotypes by genetic phasing of human families. *Am J Hum Genet*. 2011;89:382–97.
26. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The variant call format and VCFtools. *Bioinformatics*. 2011;27:2156–8.
27. Stevison LS, Woerner AE, Kidd JM, Kelley JL, Veeramah KR, McManus KF, et al. The time scale of recombination rate evolution in great apes. *Mol Biol Evol*. 2016;33:928–45.
28. O'Connell J, Gurdasani D, Delaneau O, Pirastu N, Ulivi S, Cocca M, et al. A general approach for haplotype phasing across the full spectrum of relatedness. *PLoS Genet*. 2014;10:e1004234.
29. Auton A, Fledel-Alon A, Pfeifer S, Venn O, Ségurel L, Street T, et al. A fine-scale chimpanzee genetic map from population sequencing. *Science*. 2012;336:193–8.
30. Abecasis GR, Auton A, Brooks LD, DePristo MA, Durbin RM, Handsaker RE, et al. An integrated map of genetic variation from 1,092 human genomes. *Nature*. 2012;491:56–65.
31. Gao F, Ming C, Hu W, Li H. New Software for the Fast Estimation of Population Recombination Rates (FastEPRR) in the Genomic Era. *G3 (Bethesda)*. 2016;6:1563–71.
32. Lin K, Futschik A, Li H. A fast estimate for the population recombination rate based on regression. *Genetics*. 2013;194:473–84.
33. Bansal V. Integrating read-based and population-based phasing for dense and accurate haplotyping of individual genomes. *Bioinformatics*. 2019;35:i242–8.
34. Loh PR, Danecek P, Palamara PF, Fuchsberger C, Reshef AY, Finucane KH, et al. Reference-based phasing using the Haplotype Reference Consortium panel. *Nat Genet*. 2016;48:1443–8.
35. Martin AR, Karczewski KJ, Kerminen S, Kurki MI, Sarin AP, Artomov M, et al. Haplotype sharing provides insights into fine-scale population history and disease in Finland. *Am J Hum Genet*. 2018;102:760–75.
36. Kerminen S, Havulinna AS, Hellenthal G, Martin AR, Sarin AP, Perola M, et al. Fine-scale genetic structure in Finland. *G*. 2017;7:3459–68.
37. Wang J, Santiago E, Caballero A. Prediction and estimation of effective population size. *Heredity*. 2016;117:193–206.
38. Takahata N, Satta Y, Klein J. Divergence time and population size in the lineage leading to modern humans. *Theor Popul Biol*. 1995;48:198–221.
39. Novembre J, Johnson T, Bryc K, Kutalik Z, Boyko AR, Auton A, et al. Genes mirror geography within Europe. *Nature*. 2008;456:98–101.
40. Ségurel L. The complex binding of PRDM9. *Genome Biol*. 2013;14:112.
41. Mallick S, Li H, Lipson M, Mathieson I, Gymrek M, Racimo F, et al. The Simons Genome Diversity Project: 300 genomes from 142 diverse populations. *Nature*. 2016;538:201–6.
42. Keinan A, Reich D. Human population differentiation is strongly correlated with local recombination rate. *PLoS Genet*. 2010;6:e1000886.
43. Johnston HR, Cutler DJ. Population demographic history can cause the appearance of recombination hotspots. *Am J Hum Genet*. 2012;90:774–83.
44. Kamm JA, Spence JP, Chan J, Song YS. Two-locus likelihoods under variable population size and fine-scale recombination rate estimation. *Genetics*. 2016;203:1381–99.
45. Dapper AL, Payseur BA. Effects of demographic history on the detection of recombination hotspots from linkage disequilibrium. *Mol Biol Evol*. 2018;35:335–53.
46. Spence JP, Song YS. Inference and analysis of population-specific fine-scale recombination maps across 26 diverse human populations. *Sci Adv*. 2019;5:eaaw9206.
47. Barroso VG, Puzović N, Dutheil JY. Inference of recombination maps from a single pair of genomes and its application to ancient samples. *PLoS Genet*. 2019;15:e1008449.